

(11)Publication number : 10-093655
(43)Date of publication of application : 10.04.1998

H04L 29/04
G06F 13/00

(72)Inventor : CLEMENT R ATANAASHIO
GERMAN SERGIO GOLDSCHMIDT
GURNEY DOUGLAS HOROUETIHANTO
STEPHEN EDWIN SMITH

Priority number : 96 701939 Priority date : 23.08.1996 Priority country : US

[Date of requesting appeal against examiner's decision]

Searching PAJ

of rejection]

[Date of extinction of right]

特開平10-93655

(43) 公開日 平成10年(1998) 4月10日

| | | | |
|-------------------------|-------|---------------|---------|
| (51) IntCl [*] | 識別記号 | F I | |
| H 0 4 L 29/04 | | H 0 4 L 13/00 | 3 0 3 B |
| G 0 6 F 13/00 | 3 5 5 | G 0 6 F 13/00 | 3 5 5 |

審査請求 未請求 請求項の数13 O L (全 22 頁)

(21) 出願番号 特願平9-223516

(22) 出願日 平成9年(1997) 8月20日

(31) 優先権主張番号 08/701939

(32) 優先日 1996年8月23日

(33) 優先権主張国 米国 (US)

(71) 出願人 390009531

インターナショナル・ビジネス・マシーン
ズ・コーポレーション
INTERNATIONAL BUSI
NESS MACHINES CORPO
RATION

アメリカ合衆国10504、ニューヨーク州
アーモンク (番地なし)

(72) 発明者 クレメント・アール・アタナシオ
アメリカ合衆国10568 ニューヨーク州ピ
ークスキル ディルッポ・ドライブ 5

(74) 代理人 弁理士 坂口 博 (外1名)

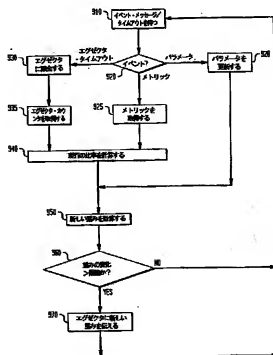
最終頁に続く

(54) 【発明の名称】 着信メッセージを経路指定する方法及びシステム

(57) 【要約】

【課題】 カプセル化クラスタの全体的なスループットを改善すること及び遠隔サービス要求の総合的な遅れを少なくすること。

【解決手段】 TCP接続ルータは各カプセル化クラスタをいくつかの仮想EC (VEC) に分割し、構成可能なポリシーにしたがって現行のサーバ負荷メトリックに基づいてVEC内に着信接続を分配することによってカプセル化クラスタリングを行う。実施の形態の1つにおいて、接続ルータはクラスタの動的な構成をサポートし、VECクライアントに中断のないサービスを与える透過的回復を可能とする。



【特許請求の範囲】

【請求項1】コンピュータ・ノードのクラスタの境界を越えて着信メッセージを経路指定する方法であって、前記クラスタが1つまたは複数のネットワークに結合されており、前記方法がポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけだし、読み取るステップと、

前記目標アドレスに基づいて、前記コンピュータ・ノードのサブセットを選択するステップと、

前記ポート番号に基づいて、前記サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択するステップと、

前記メッセージを前記経路の宛先に送るステップと、前記メッセージが境界を越えて送られている間に、前記サブセット内のメンバシップの少なくとも1つと前記サブセットの数を動的に変更するステップとを備えている方法。

【請求項2】選択が前記ポート番号と前記メッセージ・ヘッダ内のプロトコル識別子に基づいていることを備えている、請求項1に記載の方法。

【請求項3】前記サブセットの数が監視機能によって動的に変更される、請求項1に記載の方法。

【請求項4】前記サブセットのメンバシップが監視機能によって動的に変更される、請求項1に記載の方法。

【請求項5】変更が前記サブセットのメンバの置換、及び前記サブセットへのメンバの追加の少なくとも1つを含む、請求項3に記載の方法。

【請求項6】コンピュータ・ノードのクラスタを越えて着信メッセージを透過的に経路指定するシステムであって、前記クラスタが1つまたは複数のネットワークに結合されており、前記システムがポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけだし、読み取って、前記コンピュータ・ノードのサブセットを選択する手段と、

前記ポート番号に基づいて、機能を選択し、該機能が前記サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定し、該経路の宛先が前記サブセット内の前記コンピュータ・ノードである手段と、前記サブセット内のメンバシップの少なくとも1つと前記サブセットの数を動的に変更する手段とを備えているシステム。

【請求項7】前記サブセットの数が監視機能によって動的に変更される、請求項6に記載のシステム。

【請求項8】前記サブセットのメンバシップが監視機能によって動的に変更される、請求項7に記載のシステム。

【請求項9】変更が前記サブセットのメンバの置換、及び前記サブセットへのメンバの追加の少なくとも1つを含む、請求項7に記載のシステム。

【請求項10】コンピュータ・ノードのクラスタの境界

を越えて着信メッセージを経路指定する方法であって、前記クラスタが1つまたは複数のネットワークに結合されており、前記方法が境界ノードにおいて、ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読み取るステップと、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定し、該経路の宛先が前記クラスタ内のコンピュータ・ノードであるステップとを備えており、

前記境界ノードの障害を検出するステップと、障害の検出に応答して、状態情報のサブセットを前記クラスタ内の各ノードから代替境界ノードへ転送するステップと、

前記代替境界ノードにおいて、前記クラスタ内のノードから状態情報のサブセットを収集するステップと、

前記状態情報を使用して、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって送られていたのと同じ態様で、前記代替境界ノードによって分配されるようにするステップとを備えている方法。

【請求項11】ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけだし、読み取るステップと、

前記目標アドレスに基づいて、前記コンピュータ・ノードのサブセットを選択するステップとをさらに備えており、

経路の宛先が前記サブセットから選択される、請求項10に記載の方法。

【請求項12】コンピュータ・ノードのクラスタの境界ノードの障害から回復するシステムであって、前記境界ノードにおいて、ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読み取るステップと、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択し、前記経路の宛先が前記クラスタ内の前記コンピュータ・ノードである手段と、

前記境界ノードの障害を検出する手段と、障害の検出に応答して、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、前記サブセットから、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって送られていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えている代替境界ノードとを備えているシステム。

【請求項13】ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読み取り、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択し、前記経路の宛先がクラスタ内のコンピュータ・ノードである手段と、前記境界ノードの障害を検出する手段と、障害の検出に応答して、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、前記サブセットから、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって送られていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えている代替境界ノードとを備えているシステム。

【請求項14】ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読み取り、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択し、前記経路の宛先がクラスタ内のコンピュータ・ノードである手段と、前記境界ノードの障害を検出する手段と、障害の検出に応答して、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、前記サブセットから、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって送られていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えている代替境界ノードとを備えているシステム。

タ・ノードである手段と、

障害の検出に充当して、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、サブセットから、障害前の境界ノードの動作状態を再構成し、メッセージの障害前に前記境界ノードによって達成されていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えているコンピュータ・ノードのクラスターで使用される境界ノード。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明はネットワーク・コンピュータに関する。詳細に言えば、本発明はコンピュータをクラスタ化して一連の遠隔サービスをサポートすることに關する。

【0002】

【従来の技術】カプセル化クラスタ(EC)は一連のサービス(たとえば、ウェブ・サービス、NFSなど)をもたらす接続ルータ(CR)と複数のサーバ・ホストを特徴としている。カプセル化クラスタリングを提供するシステムの例には、米国特許第5371852号がある。

【0003】遠隔クライアントは、たとえばTCP/IP(たとえば、HTTP)に基づくプロトコルを使用してECにサービスを要求する。各要求に対するサービス時間はサービスのタイプ、及び対応するサーバ・アプリケーションの可用性に応じて変動する。それ故、接続の固有割振りは、利用可能なECリソースの利用度が低い偏った割振りを急速に作り出し、要求に対して不必要な遅れを導入する。

【0004】従来技術はサーバのスケールアップに關する多くの性能上の問題が存在していることを示している。たとえば、NCSAの「World Wide Web Server: Design and Performance」、IEEE Computer, 第28巻第11号、1995年11月、68-74ページ参照。ウェブ・サーバ(すなわち、httpdデモン)をサポートするのにラウンドロビンDNSを使用しているECを考えてみる。サーバはhttpを介して、ビデオ・ストリーム、データベース照会、及び静止ウェブ・ページへのアクセス・サービスを与える。各タイプの要求に対するサービス時間は、与えられるサービスのタイプ、及び關与する実際のコンテンツに応じて変動する。たとえば、複雑なデータベース照会には静止したプリロードHTMLページの何倍もの時間がかかる。要求処理時間におけるこの不釣り合いはサーバ・クラスタの偏った利用を引き起こすことがしばしばある。ラウンドロビンDNSに關連した問題は、「User Access Patterns to NCSA's Worldwide Web Server」、Kuhlth, Technical Report UIUCDSO-R-95-1394, Department of Computer Science, University of Illinois Urbana-Champaign, 1995年2月に記載されている。

【0005】従来技術はリソースの動的割振りの必要があることを示している。たとえば、「Evaluating Management Decisions via Delegation」、German Goldszmidt and Yechiam Yemini, The Third International Symposium on Integrated Network Management, San Francisco, CA, 1993年4月を参照されたい。ECは通常はホストの集合体に対して、サーバの集合体の単一のシステム・イメージを与える。しかしながら、実際のインストレーションはサービスを特定のユーザのポリシーにしたがって割り振ることを必要とすることがある。たとえば、ホストの特定のサブセットをマシナリ・ウェブ・サーバの機密トランザクションに割り振るとともに、ビデオ・オン・デマンド・サービスを、専用のハードウェアを含んでいるサーバの他のサブセットによってサポートする。

【0006】

【発明が解決しようとする課題】本発明の目的は、カプセル化クラスタの全体的なスループットを改善することである。

【0007】本発明の他の目的は、遠隔サービス要求の総合的な遅れを少なくすることにある。

【0008】本発明のさらに他の目的は、指定されたノードが障害を起こした接続ルータの動作を引き継ぎ、ネットワーク・クライアントがサービスの中断を受けないようにする手段を提供することである。

【0009】

【課題を解決するための手段】本発明の第1の態様によれば、カプセル化クラスタ(EC)はゲートウェイ・ノードとサーバ・ホストとを特徴としている。ゲートウェイ・ノードは(1)ECをいくつかの仮想VEC(VEC)に分割し、(2)構成可能なポリシーにしたがって現行のサーバの負荷に基づいて着信接続をVEC内に動的に分配し、(3)クラスタの動的な構成をサポートする。

【0010】本発明の第2の態様によれば、ゲートウェイ・ノードの障害からの迅速的な回復を可能として、クライアントへの中断のないサービスを提供するシステム及び方法が提供される。この方法によれば、クラスタまたはVEC内の各ノードはゲートウェイで保持されている状態情報のサブセットのコピーを維持する。ゲートウェイが障害を起こすと、状態情報がバックアップ・ゲートウェイへ転送される。

【0011】好ましい実施の形態において、ECは(1)VEC(すべての遠隔クライアントに対して単一のIPアドレス)、または(2)複数のVEC(いくつかのIPアドレスをエイリアス化する)として現れる。TCP-CRノードはIPアドレスを所有しており、すべてこれらのTCP接続要求を受け取る。各IPアドレスはVECと関連づけられる。TCP-CRはVECに關連づけられた重みにしたがって、新しいTCP接続

5

をホストへ分配する。TCP-CRはVECの動的定義、VECに関連づけられた重みの動的構成、VECの自動または手動管理（ホスト、サービスなどの追加、除去）を可能とする。この解決策はサーバ・ホストの動的構成、追加または除去を可能とするとともに、ネットワーク内のキャッシュされたサーバ名の問題を回避する。

【0012】

【発明の実施の形態】

1. 概要

本仮想カプセル化クラスタ・システムは米国特許第5371852号の改良形として実施できる。米国特許第5371852号は参照することにより、以下に全文が記載されているかのように、本明細書の一部となるものである。図1は米国特許第5371852号のカプセル化クラスタの発明の実施の形態を示している。米国特許第5371852号のシステムと同様に、本システムはTCP情報を送り、これはコンピュータ・クラスタの境界を越える。情報はポート・タイプ・メッセージの形態をしている。着信メッセージが送られ、サーバが応答するので、各クラスタは外部ホストに対して単一のコンピュータ・イメージとして現れる。本システムにおいて、クラスタはいくつかの仮想クラスタ（仮想カプセル化クラスタ）に分割されている。各仮想カプセル化クラスタは、クラスタ外にあるネットワーク上の他のホストに対して単一のホストとして現れる。メッセージは負荷のバランスがクラスタ・ノードのセットの間で維持されている態様で、各仮想カプセル化クラスタのメンバに送られる。

【0013】図3はTCPファミリのプロトコル用の接続ルータ、すなわちTCP接続ルータ（TCP-CR）300の実施の形態を示す。この装置は相互接続110と呼ばれる通信リンクによってまとめて接続され、クラスタを形成している2つ以上のコンピュータ・ノード（105-109）を備えている。（本発明の実施の形態の1つにおいて、相互接続はネットワークとなることに留意されたい。）ゲートウェイ109として働く、クラスタ内のコンピュータの1つはネットワーク120と呼ばれる他の通信リンクを介して、1つまたは複数の外部コンピュータもしくはクラスタ（ホスト）に接続されている。ゲートウェイを2つ以上のネットワークに接続することができ、またクラスタ内の2つ以上のノードがゲートウェイになることができる。ネットワーク、すなわち境界への各ゲートウェイの接続はネットワーク上で複数のアドレスを有することができる。各ゲートウェイはマネージャ320及びエグゼクタ340からなるTCP接続ルータ（TCP-CR）300と、図12に記載する光学式回復マネージャとを有している。マネージャはコマンド要求344をエグゼクタへ送り、応答346を評価することによって、経路を制御する。エグゼクタは米国特許第5371852号のものと同様なメッセー

6

ジ・スイッチ140及びVECルータ310からなっている。

【0014】図4は本発明の他の実施の形態を示す。好ましい実施の形態におけるのと同様に、クラスタのノード107はその応答を直接クライアント130へ返送する。しかしながら、この実施の形態においては、専用の相互接続110（図3におけるような）はなく、すべてのクラスタ・ノードが外部のネットワーク120によって接続されている。TCP接続ルータは同じものである。サンプル要求348はゲートウェイ109を介してクライアント130から、外部ネットワーク120を介してクラスタ・ノード107へ送られる。対応する応答350はノード107からクライアント130へ外部ネットワーク120を介して直接送られる。

【0015】マネージャ320の構成要素は接続割振りポリシーを実施し、仮想カプセル化クラスタの動的構成を可能とする。マネージャは動的フィードバック制御ループを介して各カプセル化クラスタのメンバにおける現在の負荷を監視し、評価する。マネージャは接続割振りポリシーを実施するものであり、このポリシーは仮想カプセル化クラスタ・サーバ全体にわたる着信接続のインテリジェントな分散を行って、クライアント要求のサービスを迅速化する。新しい重みの割当てが、クラスタ・アドミニストレータによって構成できるマネージャ・アルゴリズムによって計算される。重み割当てに対するこの判断アルゴリズムの入力は、評価された負荷メトリック、及び時間間隔などのアドミニストレータが構成できるパラメータを含んでいる。着信接続は上記の入力に基づいて各VECに動的に割り振られ、クラスタ・リソースがクライアントに最高速のサービスを与えるように割り振られるようになる。マネージャはコマンド・インタフェースも含んでおり、このインタフェースは仮想カプセル化クラスタを動的に構成するために、アドミニストレータによって使用される。マネージャの詳細については、第3節で説明する。

【0016】TCP接続ルータ109が動作を停止した場合、クラスタのすべてのノードはそれぞれの遠隔クライアントにサービスを行うことができなくなる。この問題に対処するため、機能しているゲートウェイが障害を起こした場合に、指定のバックアップ・ゲートウェイ・ノード内で活動状態となり、かつサーバ・ノードを拡張して、回復データを維持するようにする回復マネージャを追加した。クライアントがゲートウェイの障害から回復するために何らかの処置を取る必要はなく、クラスタから中断なしにサービスを受け取る。回復マネージャの詳細については、第4節で説明する。

【0017】2. エグゼクタ

図5はエグゼクタ340の好ましい実施の形態を示す。エグゼクタはコマンド・プロセッサ540、メッセージ・スイッチ140、及びVECルータ310からなる。

コマンド・プロセッサ540はエグゼクタ340に対する要求を受け取って、応答346を返す。コマンド・プロセッサはメッセージ・スイッチ140及びVECルータ310と対話を行って、要求を行い、応答を構成する。コマンド・プロセッサは接続テーブル510、VECテーブル550、ポート・テーブル520またはサーバ・テーブル530に影響を及ぼす。メッセージ・スイッチ140及び接続テーブル510は、米国特許第5371852号のメッセージ・スイッチ及び接続テーブルと同じものである。本発明の好ましい実施の形態において、VECルータ310は着信パケットを修正しない。パケットはサーバへ送られるが、サーバは応答が内部ノードからクライアントへ直接送られるように構成されている。

【0018】メッセージ・スイッチ140は本質的に、米国特許第5371852号のメッセージ・スイッチと同じものである。しかしながら、本発明のため、好ましい実施の形態のメッセージ・スイッチは最適化され、付加的なチェックがメッセージ・スイッチに追加されている。メッセージ・スイッチはメッセージがVECルータに既知のVECに対するものであるかどうかを調べるためにチェックを行わなければならない。

【0019】VECルータは外部ネットワーク上のクライアントに対する各VECを表す一連のアドレスを維持している。VECルータは要求をクラスタの内部ノードへ送るが、受け取った要求を修正しない。クラスタの各内部ノードは1つまたは複数のVECに関連づけられており、関連づけられているVECに対する要求だけを受け取る。当分野で周知の技法を使用して、本発明においては、VECを表すアドレスに送られたパケットを受け入れ、クライアントに直接応答するように、内部ノードを構成する。従来技術においては、メッセージ・スイッチ140は着信要求（図1の140）に合わせてパケット・ヘッダを書き換え、要求に対する応答（図1の120）に合わせてパケット・ヘッダを書き換えなければならない。本発明においては、パケット・ヘッダの書き換えは必要ない。（従来技術を本発明とともに使用することはできる。）パケット・ヘッダが書き換えられず、応答パケットがゲートウェイ・ノード109を通して流れないため、本発明の性能は従来技術よりも良好なものである。応答パケットがTCP接続ルータを通して流れないため、メッセージ・スイッチはクラスタ内部のノードから応答パケットを受け取ることはない。その結果、好ましい実施の形態において、ヘッダの書き換えが排除され、内部ノードからの応答パケットに対するチェックが排除される。

【0020】この改善の直接的な結果として、VECが調べるのはクライアントと、サービスを与える内部ノードとの間の流れの半分だけとなる。これは正確な接続テーブルを維持することを困難とする。この問題を解決す

るために、本発明はその接続テーブルに固有な2つの新しいタイム、状態タイムアウトとFINタイムアウトを使用する。これら2つのタイムと通信流、及び当分野で周知のタイムを使用すると、接続テーブルを正確に維持することができる。

【0021】接続テーブル項目は2つの状態ACTIVE及びFINの1つであると考えられる。新しい接続が確立された場合には、接続テーブル項目が作成され、活動状態にされる。パケットが接続テーブル内に項目がある接続を流れている場合、接続項目にはタイム・スタンプがつけられる。VECルータがクライアントからサービスを与えるノードへのFIN流を見つけた場合、関連する接続テーブル項目がFIN状態にされる。（パケットはFIN状態にされた接続上を流れることができる。）

接続テーブル項目は閉じられており、最後のパケットがクライアントからその接続上のサーバに送られてから、FINタイム・アウトによって特定されている時間数が経過した場合に、バージを行うのに利用できると考えられる。クライアントがFINを送出することなく障害を起こした場合、接続レコード項目が残る。失効したタイムアウトは最後のパケットが活動会話を流れてからどれくらい待ってから、接続テーブル項目をバージするかを指定する。

【0022】図7-図9はVECルータ310の流れ図を示す。図7において、VECルータはパケット702を待っている。パケットを受け取ると、VECルータはパケットが既存のTCP接続に対するものであるか、新しいTCP接続に対するものであるかをチェックする（740）。パケットが既存のTCP接続に対するものである場合、パケットがFINであるか、SYNであるか、RSTであるか（すべてのパケット・タイプは当分野で周知のものである）をチェックする（708）。パケットがこれらのものの1つでない場合には、パケットを接続に関連づけられた内部ノードに送る（722）。これ以外の場合には、パケットがRSTであるかどうかをチェックする（701）。パケットがRSTである場合には、会話が接続テーブルからバージされて、接続をリセットし（712）、パケットはその接続と関連づけられている内部ノードに送られる（722）。パケットがRSTでない場合には、VECルータ310はパケットがSYNであるかどうかをチェックする（714）。パケットがSYNである場合には、接続を確立し（716）、この接続が以前に存在していない、その接続を活動状態にする。VECルータ310は次いで、パケットがFINであるかどうかをチェックする（718）。パケットがFINである場合には、接続はFIN状態となる（720）。FIN処理後、あるいはパケットがFINでない場合、パケットは接続に関連づけられているサーバへ送られる（722）。

【0023】図8は非存在接続の流れ図を示す。チェッ

ク704で非存在接続が見つかった場合、VECルータはまずパケットがSYNであるかどうかをチェックする(724)。パケットがSYNでなければ、廃棄する(726)。パケットがSYNである場合、接続が活動状態にセットアップされ(728)、サーバが選択される(730)、パケットが選択されたサーバへ送られる(722)。

【0024】図9は新しい接続に対してサーバを選択する(730)処理の流れ図を示す。本発明において、この機能は加重経路指定を実現する。サーバの選択に関するこの検討のために、VECの内部ノードには1からnまでの番号がつけられているものとする。たとえば、VECに7つのノードがある場合、番号は1、2、3、4、5、6及び7となる。サーバの選択に関するこの検討のために、適格な重みを適正な最大値から1までの番号がつけられているものとする。たとえば、最大適正値が5である場合、適格な重みは5、4、3、2、及び1となる。0は特別な値である。重みは降順で選択することもできる。本発明は重みを個々のサービスを与える各内部ノードに関連づける。これは各サービスに対して、少なくとも1つのノードが最大の非ゼロの重みを有しているか、あるいはすべてのノードがゼロの重みを有していることを保証する。

【0025】サーバを選択する機能(730)はまず、次にもっとも高いサーバ(734)と現行の適格な重みを拾い出す。次いで、この値が大きすぎないものかどうかをチェックする(735)。この値が大きすぎない場合には、この値に対応するサーバが良好な選択であるかどうかをチェックする(746)。(このチェックについては、後述する。)この値が大きすぎる場合には、最初のサーバ(736)と次に低い重みを拾い出す。次いで、次に低い重みがゼロであるかどうかをチェックする(738)。次に低い重みがゼロでない場合には、現行の適格な重みの代わりに、これを使用する。この機能は現行のサーバが良好な選択であるかどうかをチェックする(746)。最初のサーバと最大の重みを選択した後、この機能はパケットを送るのに利用できるサーバがあるかどうかをチェックする(742)。利用可能なすべてのノードの重みがゼロである場合、サーバはまったく利用できない。利用できるサーバがない場合には、サーバを選択することなく、パケットは戻される(744)。利用可能なサーバがある場合には、この機能は良好な選択があるかどうかをチェックする(746)。良好な選択は、重みが現行の適格な重みより等しいか、それよりも大きいサーバとして定義されている。これが良好な選択である場合には、サーバが選択され(748)、VECルータへ展される(750)。これが良好な選択でない場合には、アルゴリズムは次のサーバを拾い出す(734)。

【0026】最大の重みが非ゼロであり、少なくとも1

つのノードが最大の重みを有しているか、あるいはすべてのノードの重みがゼロであるため、サーバ選択機能は常に終了する。正の重みのあるノードがある場合、サーバを選択する機能は重みの比に基づいてパケットを分配する。たとえば、任意の2つのノードの間で、一方の重みが3であり、他方の重みが2である場合、重みが3のノードは、重みが2のノードが受け取る2つのパケットごとに3つのパケットを取得する。

【0027】図10はVECルータが使用するデータ構造の実施の形態を示す。VECテーブル550は外部ネットワーク上のVECアドレスである一連のアドレスを含んでいる。VECに個々に関連づけられるすべてのパラメータも、このテーブルに含まれている。各VECテーブルはポート・テーブル520に関連づけられており、このポート・テーブルはVECがサービスを与えている一連のポート802を含んでいる。各ポート項目802には失効タイムアウト804、FINタイムアウト806及びその他のポート固有の属性808が関連づけられている。各ポートには、このポートに関連づけられているサービスを与えるために使用されるVECの内部ノードのサブセットが関連づけられている。ノード・テーブル530はポートに関連づけられているノード820のアドレス、このノードに関連づけられている現行の重み822、及びその他のノード固有の情報830を含んでいる。(ノード固有の情報の例は、活動状態の接続の数、FIN状態の接続の数、及び完成した接続の総数を示すカウンタである。)ノード・テーブル530はこのテーブル内のノードのセットに対して加重経路指定を実施するためにサーバを選択する機能に必要な状態も含んでいる。ノード・テーブルはノード810の総数、最後に選択されたノード812、現行の適格な重み814、最大重み816、及び重みのバンド818を含んでいる。重みのバンドは最大重みの変動を制限するために使用される。いずれのノードも、重みのバンドよりも大きい重みを持つことはできない。

【0028】3. マネージャ
接続ルータマネージャ(マネージャ320)の発明は構成可能なポリシーにしたがって、いくつかの負荷メトリックを使用して着信接続を動的に分配する方法及び装置である。マネージャはエグゼクタ340の経路指定アルゴリズムの重みを動的に修正して、クラスタ・リソースの割振り最適化を制御する制御ループを備えている。本発明の目的は、クラスタの現行状態にしたがって着信TCPを分配することによって、クラスタの全体的なスループットを改善し、サービス要求の総合的な遅れを少なくすることである。それ故、本発明は接続をサーバ・ホストへ分配する、サーバの利用度を改善し、要求のサービスを行うことの遅れを少なくする方法を記載する。

【0029】図6は5つのノード(105、106、107、108、及び109)のクラスタ600内の、本

発明のマネージャ320の実施の形態のサンプルを示す。図6は図4の代替ネットワーク構成を使用しているが、図3の構成も可能である。ノード1つはゲートウェイ109であり、これは外部ネットワーク120に接続しており、TCP接続ルータ300（エグゼクタ340及びマネージャ320）を実行する。マネージャ320は5つの一般的な構成要素、負荷マネージャ（Muddy）610、外部制御インタフェース（Callbuddy）620、クラスタ・ホスト・メトリック・マネージャ（Hostmonitor）630、順方向メトリック・ジェネレータ（FMG）640、及びユーザ・プログラマブル・メトリック・マネージャ（UPMM）650からなっている。

【0030】Muddy610は入力メトリック、ホスト・メトリック、サービス・メトリック、及びユーザ・メトリックという4つの異なるクラスのメトリックを使用して、エグゼクタ10の重み機能を実行することができる。Muddy610はこれらのメトリックとその他の関連する情報を、エグゼクタ・インタフェース346、Callbuddyインタフェース624、Hostmonitorインタフェース634、FMGインタフェース644、及びUPMMインタフェース654から受け取る。Muddyはインタフェース344を介して、各VECポート・サーバに対するエグゼクタ経路指定アルゴリズムに関連づけられた重みを制御する。

【0031】Muddy610は定期的にエグゼクタ340に、インタフェース346を介して、各サーバに関連づけられた内部カウンタの値を要求する。たとえば、各サーバに対して確立された接続の総数に関するカウンタの値を定期的に要求する。時間T1及びT2でポーリングされた2つのカウンタを減算することにより、Muddy610は時間間隔T1-T2中に受け取った接続の数を表すメトリック変数を計算する。このような入力メトリックの集合は各VEC及び各ポート・サーバに対する接続要求の特性速度の概数を与える。

【0032】Hostmonitor630は定期的にMuddy610に、メッセージ・インタフェース634を介して、クラスタ内の各ホストの状態に関する情報を送る。この状態情報を得るには、さまざまな周知の方法がある。たとえば、Hostmonitorはプログラム・スクリプトを実行して、ホスト固有のメトリックを評価する監視エージェント635を使用することができる。たとえば、スクリプトはネットワーク接続に対するメモリ・バッファの現行レベルの利用度を評価することができる。メトリック・レポートをポリシー固有の閾値時間内に受け取らなかった場合には、対応するホスト・メトリックに特別な値が与えられ、マネージャはホストに到達できないと判断することができ、それ故、それ以上の接続要求はこのホストに送られない。Hostmonitor630はすべての監視エージェントのレ

ポートを調整し、これをMuddyに提示する。

【0033】順方向メトリック・ジェネレータ（FMG）640は順方向要求、すなわちゲートウェイ109のコンピュータから始まる要求を使用して、アプリケーション固有のメトリックまたはサービス固有のメトリックを作成し、評価する。評価はクラスタ・ホスト・サーバの各々に対して適切な要求を作成し、これらの応答の遅れを測定することからなっている。たとえば、HTTPサーバにおける順方向遅延メトリックを取得するために、FMGは特定のポート（たとえば、ポート80）にサービスを行うクラスタにおける各HTTPサーバに対してHTTP「GET /」を生成する。次いで、FMG640はHTTP要求にサービスを行う際の対応する遅れを測定し、メトリックス・ベクトルをMuddy610へ送る。ポリシー固有の閾値時間までに要求に回答がない場合には、FMGは対応するサービス・ノードに、特定のサービス・タイプの新しい要求を一時的に受け取らないとのマークをつける。この情報をマネージャが使用して、特定のホストにおけるサービスに一時的に到達できないと判断し、それ故、このタイプのこれ以外の接続要求はこのサービスに送られない。

【0034】ユーザ・プログラマブル・メトリック・マネージャ（UPMM）650により、本発明のユーザが任意の新しいメトリックを接続の管理のために考えることを定義できるようにする。このようなメトリックは所与のクラスタ・インストールが実施を望むことのできる任意のポリシーを記述できる。たとえば、任意のポリシーは管理上の考慮事項のため、ある期間中にあるセットのクラスタ・ホストがTCP接続を受け取ってはならないことを必要とすることができる。UPMM650はインタフェース654を介して、メトリックとしてこれらのポリシーをMuddyに通信する。

【0035】Callbuddy620の構成要素はアドミニストレータが、Muddy610のパラメータの任意のものを動的に調節することを可能とする。Callbuddyにより、アドミニストレータがアルゴリズムを構成して、Muddyによって実施された重みの割当てを計算するようにすることが可能となる。たとえば、アドミニストレータは現行メトリックの各々に関連づけられた重みを動的に変更しようとすることができる。アドミニストレータは、たとえば、（1）ホスト・メトリックの重みを高くし、（2）サービス・メトリックの重みを低くし、（3）入力メトリックに合わせてエグゼクタ340の重みをポーリングする周波数を高くすることを選択できる。Callbuddy620の構成要素はインタフェース622を介してアドミニストレータの要求を受け取り、インタフェース624を介してMuddy610に通知する。

【0036】Muddy610の構成要素は負荷マネージャであり、サーバと接続ルータ・ゲートウェイ・ノ

ードと間の動的フィードバック制御ループを確立する。Mbuddyはエグゼクタ610の経路指定アルゴリズムの重みを調節するので、負荷メトリックにしたがって軽負荷がかけられているサーバは、そのタイプの着信TCP接続の大きい部分を受け取る。上記で定義されたような任意のセットの負荷とポリシー・メトリックが与えられると、Mbuddyはその現行のメトリック及びその現行の重みに基づいて、各VEC内の各ポートの各サーバに対する新しい相対重みを計算する。

【0037】あらゆるVECにおける各ポートに対する重みの割当ては次のようにして計算される。(1)すべての実行サーバに対するすべての総合メトリック(AM)を計算する。(2)各実行サーバ(CWP)に対するすべての現行の重みの比率を計算する。(3)各メトリックMに関し、各サーバSについてその値のメトリック比率(MP)を計算する(総合AMに関して)。

(4)各サーバに関し、新しい重みNWを計算する。
(4a)サーバが静止している場合には、NWを0にセットする。(4b)サーバがスティッキ重みWを有している場合には、Wの値をNWとして使用する。(4c)次の式によって、ベクトルNWVを計算する(ただし、各項目NWV[i]は単一のメトリックM[i]に基づく)。

$NWV[i] = AW + [(CWP - MP) / SF]$
ただし、AWは重みの現行範囲内での平均重みであり、SFは平滑化計数のパラメータである。(4d)各サーバの新しい重みNWを次のものとして計算する。

$NW = NWV[1] * W[1] + NWV[2] * W[2] + \dots + NWV[i] * W[i]$

【0038】図11は本発明により、メトリックがどのようにメトリックによって受け取られるのか、また重みの割当てがどのように計算されるのかについての流れ図の説明である。最上部の枠910はマネージャMbuddyがメッセージまたはタイムアウトのいずれかであるイベントを待機していることを示す。判断ブロック920は発生するイベントのタイプを決する。値のリフレッシュを必要とするタイムアウトがあった場合、ブロック930において、エグゼクタに照会を行って、入力メトリックを与える一連のカウント値を取得する(935)。ブロック920で、イベントがパラメータの更新についての要求であると判断された場合、対応するパラメータが更新される(928)。たとえば、アドミニストラータはメトリックに関連づけられた重み、またはポーリング期間を更新することができる。ブロック920で、イベントがメトリック更新の受取りであったと判断された場合には、ブロック925において、メトリックが検索され、これにしたがって、内部変数値がセットされる。新しいメトリックが到着した場合には、ブロック940において、アルゴリズムがすべてのメトリックの現行の比率と、現行の重みを計算する。次いで、ブロック

950において、サーバ・ノードの各々に対する新しい重みNWが、上述の式を使用して計算される。このブロックは各サーバiが重み項目を有している重みNW[i]の新しいベクトルを作成する。判断ブロック960は任意の閾値機能によって、重みNW[i]の計算された新しいベクトルが、現行の重みベクトルと異なっているかどうかを判断する。新しいベクトルが異なっている場合には、ブロック970において、エグゼクタに新しい重みが通知され、それ以外であれば、アルゴリズムはTOP状態に戻り、新しいイベントを待機する。

【0039】4. 回復マネージャ
機能しているゲートウェイの障害を検出すると、指定されたバックアップ・ゲートウェイにおける回復マネージャが活動状態となる。障害検出は従来、A. Bhidhe, 「A highly Available Network File Server」、USENIX Conference, 1991年冬期、テキサス州ダラス、第199ページ、またはF. Jahaniand, 「Processor Group Membership Protocols: Specification, Design and Implementation」、Proceedings of the 12th Symposium of Reliable Distributed Systems, 第2-11ページ、Princeton, ニュージャージー州、1993年10月、IEEE Computer Societyに記載されているように行うことができる。

【0040】回復マネージャはまずHA/NFSで教示されているようにして、障害を起こしたゲートウェイからネットワーク接続を除去し、次いで、すべての活動サーバ・ノードに問い合わせを行って、それぞれのシャドウ接続テーブルから状態情報を取得し、この情報からゲートウェイのメッセージ・スイッチ内に接続テーブルを構成する。引継プロセスはTCP/IPのタイムアウト期間内に完了しなければならない。その結果、既存の接続が失われることはない。これを達成するために、内部ノードは新規なハイブリッド・アルゴリズム(後述)を実行して、接続が非活動状態になった時点を感じ、これをそれぞれのシャドウ接続テーブルから除去し、その結果、活動接続が引き継ぎゲートウェイに対して記述される。すべての機能中のクラスター・ノードが応答すると(指定の期間内に応答しないノードは機能していないものと見なされる)、バックアップ・ゲートウェイで実行されている回復マネージャはそれ自体のネットワーク・インタフェースを使用可能とし、クラスターのIPアドレス宛のパケットを受け取る。この最後のステップはバックアップ・ゲートウェイが動作するのを可能とするのに必要な作業を完了する。マネージャの構成要素によって使用される相対的に静的な構成データは、主ゲートウェイバックアップ・ゲートウェイの間で共用されているファイルに維持され、引継中にバックアップによって読み取られる。

【0041】自明ではあるが、受け入れられない解決策は接続情報をバックアップ・ゲートウェイに重複して維

持することである。これは各確立された接続及び各終了した接続における主ゲートウェイバックアップ・ゲートウェイの間の「2段階」プロトコルを必要とし、パフォーマンスのコストが厳しいものであるため、拒絶された。

【0042】図12は高可用性のゲートウェイを備えたカプセル化クラスタの構成を示す。主ゲートウェイ1050は外部ネットワーク120に活動状態で接続されている。指定されたバックアップ・ゲートウェイ1030はネットワーク120に対する物理的ではあるが、非活動状態の接続を含んでいる。正規のカプセル化クラスター・ゲートウェイの構成要素、マネージャ320及びエグゼクタ340に加えて、各ゲートウェイは回復マネージャ1020を含んでいる。(主ゲートウェイは障害及び回復後にバックアップとなることができる。)各サーバ・ノード107はシャドウ接続テーブル1010を含んでおり、このテーブルには、外部ネットワーク120に対するその活動接続に関する情報が維持されている。

【0043】メッセージ(ipパケット)がクラスター・ゲートウェイに到着し、特定のTCPまたはUDPプロトコル・ポートへ送られる。ゲートウェイ内のメッセージ・スイッチにより、メッセージ経路指定機能をプロトコル・ポートに合わせてインストールすることが可能となる。経路指定機能が関連づけられたポートに到着する各メッセージに対して呼び出され、内部ノード及びメッセージが送られるポートの選択を担当する。確立済みの接続を指定する情報及び接続を保持するクラスタが、ゲートウェイ内のテーブルに記録される。このテーブルをメッセージ・スイッチが使用して、確立済みの接続上の着信パケットを適正なクラスター・ノードへ送る。

【0044】どのサーバ・ポートがインストールされたメッセージ・スイッチ機能を持っているかなどの相対的に静的な情報が維持され、他のマネージャ構成情報が、主ゲートウェイ及びバックアップ・ゲートウェイ両方にアクセス可能な共用ファイルに保管される。現行の接続情報はきわめて迅速に変化し、本明細書に記載する技法にしたがって管理される。

【0045】各内部ノード107はそれ自体の接続のために(他のノードに対する接続に対してではなく)ゲートウェイ経路指定テーブルのシャドウ1010を維持している。このシャドウ・テーブルをノードが使用して、引継中にバックアップ・ゲートウェイ1030における回復マネージャ1020からの、引継ゲートウェイの要求に応答する。このテーブルは内部ノードが引継ゲートウェイに応答するのに必要とする時間を大幅に少なくし、これは確立済みの接続を活動させておくためには、接続ベースのプロトコルが通信を正常に完了できるようにする「タイムアウト」期間内に引継ゲートウェイが作動していなければならないため、きわめて重要である。

【0046】接続テーブルの項目に対するスペースを再

利用するため、ゲートウェイ内、及び内部ノードに維持されているシャドウ内の両方で、次のように進める。接続は2つの状態、すなわちACTIVEまたはFINのうちのいずれかになっている。接続テーブル項目には参照ごとにタイム・スタンプがつけられる。FIN_TIME_OUTと呼ばれるユーザ構成可能なタイムが維持されている。このタイムは閉鎖されると考えられるFIN状態における会話に対する最後の参照後の時点を表している。タイムはサービス・アドレスごと、あるいはポートごとのグローバルなものでよい。活動クローズ(接続の一方の側がFINを送ったが、他方が接続上で送信を継続している)の意図は、サーバがクライアントへのデータの送信を継続することができ、データ電送の終了時に、会話が閉鎖されることである。クライアントはサーバにクライアントからこれ以上要求が送られないということをサーバに伝える手段として会話を能動的に閉鎖する。ここでの検討のために、クライアントの要求がルータを経由して送られているものと想定する。このプロトコルが機能するのは、サーバが肯定応答を受けるデータを送り続けているからである。ルータは、したがってサーバは肯定応答を見て、接続テーブル項目に連続的にタイム・スタンプをつける。サーバがクライアントに対するデータの送信を完了し、その会話の「半分」を閉鎖すると、最後の肯定応答がクライアントからサーバへ向かって流れる。FIN_TIME_OUTが経過した後、サーバは接続項目をバージングすることができる。第2のタイムSTALE_TIME_OUTがゲートウェイによって維持されている。STALE_TIME_OUT以外には何の活動もしていない活動状態の何らかの接続はバージングできる。

【0047】このアルゴリズム(接続再構成アルゴリズム)は、バックアップゲートウェイによる引継プロセスのサポートを内部ノードが維持している接続テーブルのシャドウ内の項目に対するスペースを再利用するために、内部ノードにおいても実行される。このようにして、シャドウ・テーブル内の項目の数をできるだけ少ない数に維持し、これにより引き継ぎプロセスをできるだけ迅速に進めることができる。

【0048】デフォルトとして、FIN_TIME_OUTをTCPの最小セグメント超(MSL)の3倍の値にセットする必要がある。デフォルトのSTALE_TIME_OUTはTCPの失効タイムアウトよりも長くなければならない。タイムが関連づけられているプロトコルを考慮することにより、FIN_TIME_OUTに対するより妥当な値を見つけたことができる。

【0049】主ゲートウェイが障害を起したと判断されたか、あるいは明示のアドミニストレーション上の判断によって判断されたかしたため、バックアップ・ゲートウェイが活動状態とならなければならないと判断された場合、以下のステップがバックアップ・ゲートウェイ

17

イ1030内の回復マネージャ1020によって取られる。

【0050】(1) i p アドレス引継を使用し、バックアップ・ゲートウェイは主ゲートウェイのネットワーク接続を除去する。ダウンしたと想定されるゲートウェイが実際に、ネットワークからのメッセージを受け入れないようにするため、このステップが必要である。このステップを行わないと、ある種の障害、すなわち「部分的に障害を起こしたゲートウェイ」がメッセージの受取りを継続し、処理を行うという障害が生じる可能性がある。これはシステムの安全性を危険にさらすのである。

【0051】(2) バックアップ・ゲートウェイがクラスタの機能している各ノードに問い合わせを行い、それぞれノードに割り振られているすべてのUDPポート、主ゲートウェイクラスタ外のホストとの間の主ゲートウェイを介して確立されたTCP接続の記述を要求する。バックアップ・ゲートウェイは私用i p ベース・プロトコルを使用してこれを行う。各ノードに維持されているシャド接続テーブルはノードからの即時応答を可能とし、確立済み接続がゲートウェイの引継時にタイムアウトとならない確率を高くする。閉鎖した接続を認識し、これらをサポートするために使用されるスペースを再利用するための上述したアルゴリズムは、シャド接続テーブルのサイズを最小限のものとし、ゲートウェイの引継を達成するのに必要な時間を短縮するのに寄与する。

【0052】(3) バックアップ・ゲートウェイが機能している各ノードからの応答を記録し、ノードのUDPポート及びTCP接続を、図12のバックアップ・ゲートウェイのエグゼクタ340内の図5の接続テーブル510に記録する。

【0053】(4) すべての機能しているクラスタ・ノードが応答した場合(指定された期間内に応答しないノードは機能していないものと見なされる)、バックアップ・ゲートウェイはそれ自体のネットワーク・インタフェースを使用可能として、クラスタのi p アドレス宛のパケットを受け取るようにする。この最後のステップはバックアップ・ゲートウェイが作動状態となることを可能とするのに必要な作業を完了する。

【0054】本発明を好ましい実施の形態により説明してきたが、各種の改変及び改善が当分野の技術者には思い浮かぶであろう。それ故、好ましい実施の形態が例として挙げられたものであって、限定事項として挙げられたものではないことを理解すべきである。本発明の範囲は首記の特許請求の範囲によって画定されるものである。

【0055】まとめとして、本発明の構成に関して以下の事項を開示する。

【0056】(1) コンピュータ・ノードのクラスタの

18

境界を越えて着信メッセージを経路指定する方法であって、前記クラスタが1つまたは複数のネットワークに結合されており、前記方法がポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけたし、読み取るステップと、前記目標アドレスに基づいて、前記コンピュータ・ノードのサブセットを選択するステップと、前記ポート番号に基づいて、前記サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択するステップと、前記メッセージを前記経路の宛先に送るステップと、前記メッセージが境界を越えて送られている間に、前記サブセット内のメンバシップの少なくとも1つと前記サブセットの数を動的に変更するステップとを備えている方法。

(2) 選択が前記ポート番号と前記メッセージ・ヘッダ内のプロトコル識別子に基づいていることを備えている、上記(1)に記載の方法。

(3) 前記サブセットの数が監視機能によって動的に変更される、上記(1)に記載の方法。

(4) 前記サブセットのメンバシップが監視機能によって動的に変更される、上記(1)に記載の方法。

(5) 変更が前記サブセットのメンバの置換、及び前記サブセットへのメンバの追加の少なくとも1つを含む、上記(3)に記載の方法。

(6) コンピュータ・ノードのクラスタを越えて着信メッセージを透過的に経路指定するシステムであって、前記クラスタが1つまたは複数のネットワークに結合されており、前記システムがポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけたし、読み取り、前記コンピュータ・ノードのサブセットを選択する手段と、前記ポート番号に基づいて、機能を選択し、該機能が前記サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定し、該経路の宛先が前記サブセット内の前記コンピュータ・ノードである手段と、前記サブセット内のメンバシップの少なくとも1つと前記サブセットの数を動的に変更する手段とを備えているシステム。

(7) 前記サブセットの数が監視機能によって動的に変更される、上記(6)に記載の方法。

(8) 前記サブセットのメンバシップが監視機能によって動的に変更される、上記(7)に記載の方法。

(9) 変更が前記サブセットのメンバの置換、及び前記サブセットへのメンバの追加の少なくとも1つを含む、上記(7)に記載の方法。

(10) コンピュータ・ノードのクラスタの境界を越えて着信メッセージを経路指定する方法であって、前記クラスタが1つまたは複数のネットワークに結合されており、前記方法が境界ノードにおいて、ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけたし、読み取るステップと、前記ポート番号に基づいて

50

て、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定し、該経路の宛先が前記クラスタ内のコンピュータ・ノードであるステップとを備えており、前記境界ノードの障害を検出するステップと、障害の検出に応じて、状態情報のサブセットを前記クラスタ内の各ノードから代替境界ノードへ転送するステップと、前記代替境界ノードにおいて、前記クラスタ内のノードから状態情報のサブセットを収集するステップと、前記状態情報を使用して、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって達成されていたのと同じ態様で、前記代替境界ノードによって分配されるようにするステップとを備えている方法。

(11) ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号及び目標アドレスを見つけだし、読み取るステップと、前記目標アドレスに基づいて、前記コンピュータ・ノードのサブセットを選択するステップとをさらに備えており、経路の宛先が前記サブセットから選択される、上記(10)に記載の方法。

(12) コンピュータ・ノードのクラスタの境界ノードの障害から回復するシステムであって、前記境界ノードにおいて、ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読取り、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択し、前記経路の宛先が前記クラスタ内の前記コンピュータ・ノードである手段と、前記境界ノードの障害を検出する手段と、障害の検出に応じて、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、前記サブセットから、障害前の前記境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって達成されていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えている代替境界ノードとを備えているシステム。

(13) ポート・タイプ・メッセージのメッセージ・ヘッダ内のポート番号を見つけだし、読取り、前記ポート番号に基づいて、サブセット内の複数の考えられる宛先からメッセージに対する経路の宛先を判定する機能を選択し、前記経路の宛先がクラスタ内のコンピュータ・ノ

ードである手段と、障害の検出に応じて、前記クラスタ内の各ノードから状態情報のサブセットを収集する手段と、サブセットから、障害前の境界ノードの作動状態を再構成し、メッセージが障害前に前記境界ノードによって達成されていたのと同じ態様で、代替ノードによって分配されるようにする手段とを備えているコンピュータ・ノードのクラスタで使用される境界ノード。

【図面の簡単な説明】

【図1】従来技術のカプセル化クラスタ・システムを示す図である。

【図2】従来技術のメッセージ・スイッチを示す図である。

【図3】本発明の実施の形態による仮想カプセル化クラスタ・システムを示す図である。

【図4】本発明の他の実施の形態による仮想カプセル化クラスタ・システムを示す図である。

【図5】図3及び図4のエグゼクタの詳細な図である。

【図6】図3及び図4のマネージャの詳細な図である。

【図7】エグゼクタの流れ図である。

【図8】エグゼクタの流れ図である。

【図9】エグゼクタの流れ図である。

【図10】エグゼクタのデータ構造を示す図である。

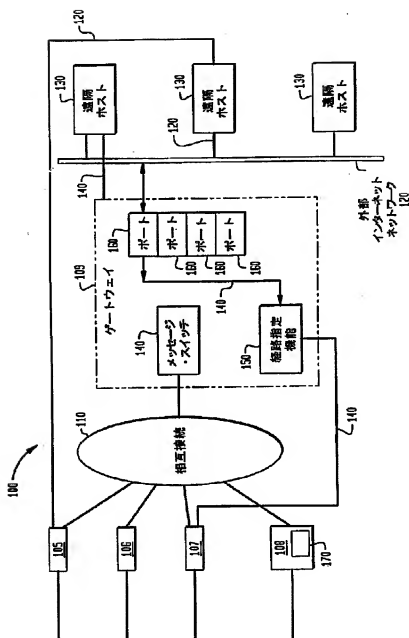
【図11】マネージャの流れ図である。

【図12】本発明の実施の形態による高可用性ゲートウェイを有するクラスタを示す図である。

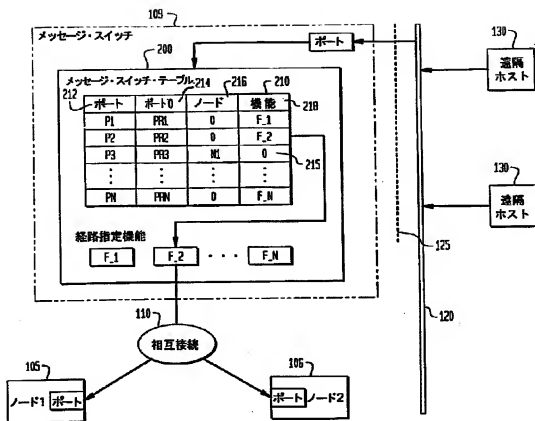
【符号の説明】

- 105 コンピュータ・ノード
- 106 コンピュータ・ノード
- 107 コンピュータ・ノード
- 108 コンピュータ・ノード
- 109 ゲートウェイ
- 110 相互接続
- 120 ネットワーク
- 140 メッセージ・スイッチ
- 300 TCP接続ルータ(TCP-CR)
- 310 VECルータ
- 320 マネージャ
- 340 エグゼクタ

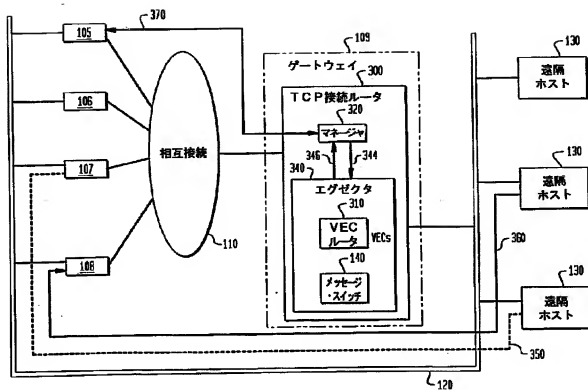
【図1】



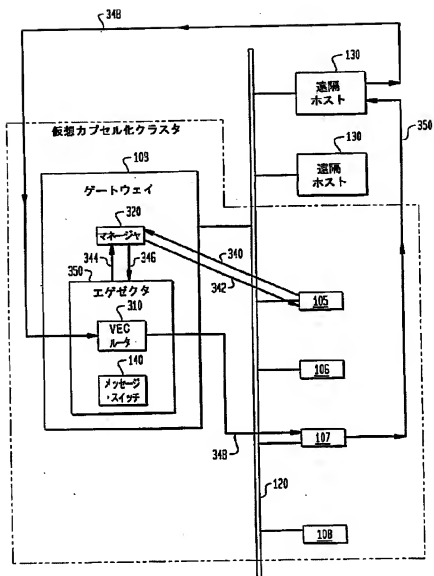
【図2】



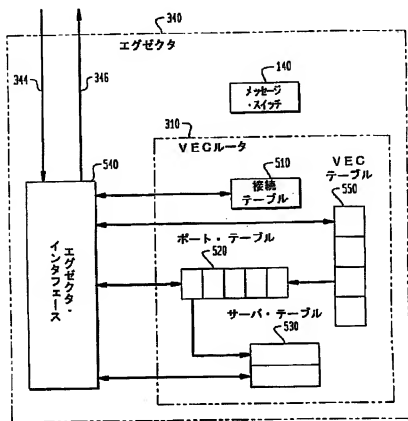
【図3】



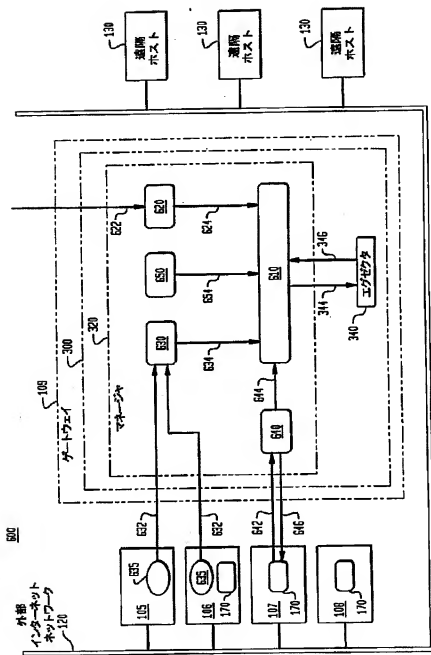
【図4】



【図5】

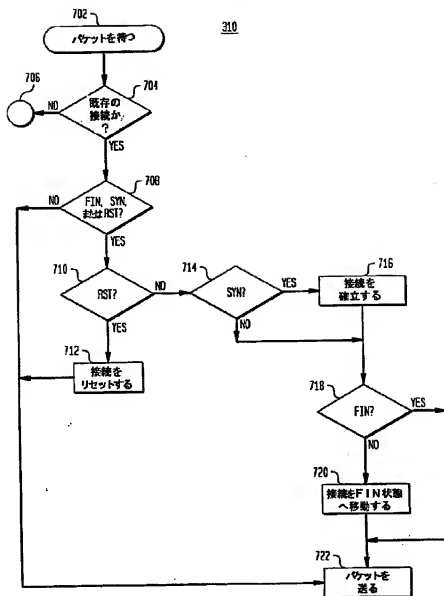


【図6】

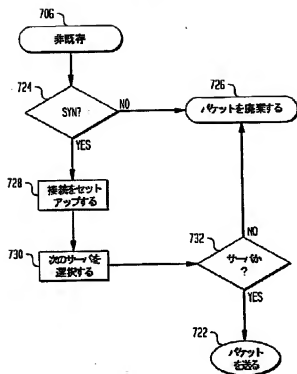


【図7】

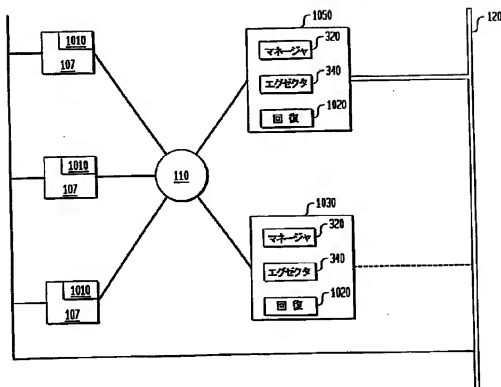
310



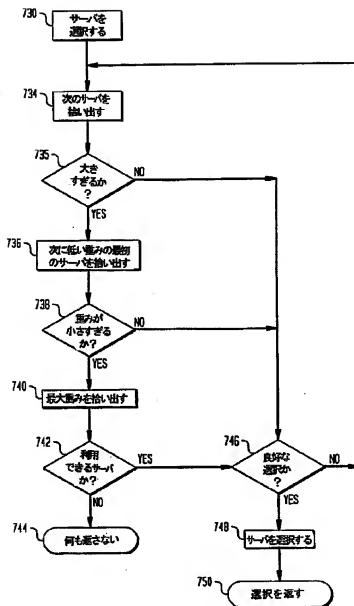
【図8】



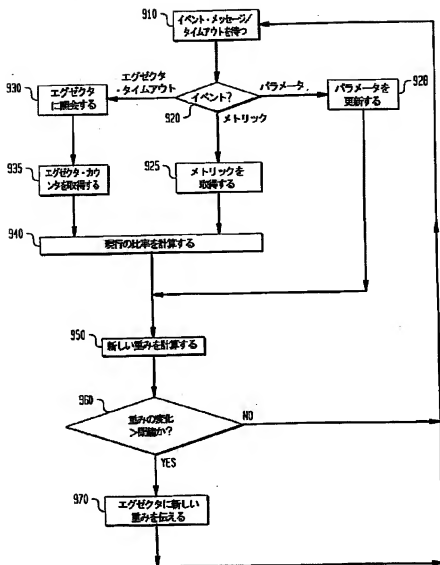
【図12】



【図9】



【図11】



フロントページの続き

(72)発明者 ジャーマン・セルジオ・ゴールドシュミット
アメリカ合衆国10522 ニューヨーク州ド
ッブス・フェリー チェストナット・リッ
ジウェイ 21

(72)発明者 ガーニイ・ダグラス・ホロウェイ・ハント
アメリカ合衆国12590 ニューヨーク州ワ
ッピンジャー・フォールス エッジヒル・
ドライブ 127

(72)発明者 ステファン・エドウィン・スミス
アメリカ合衆国10541 ニューヨーク州マ
ホバックハットフィールド・ロード 19